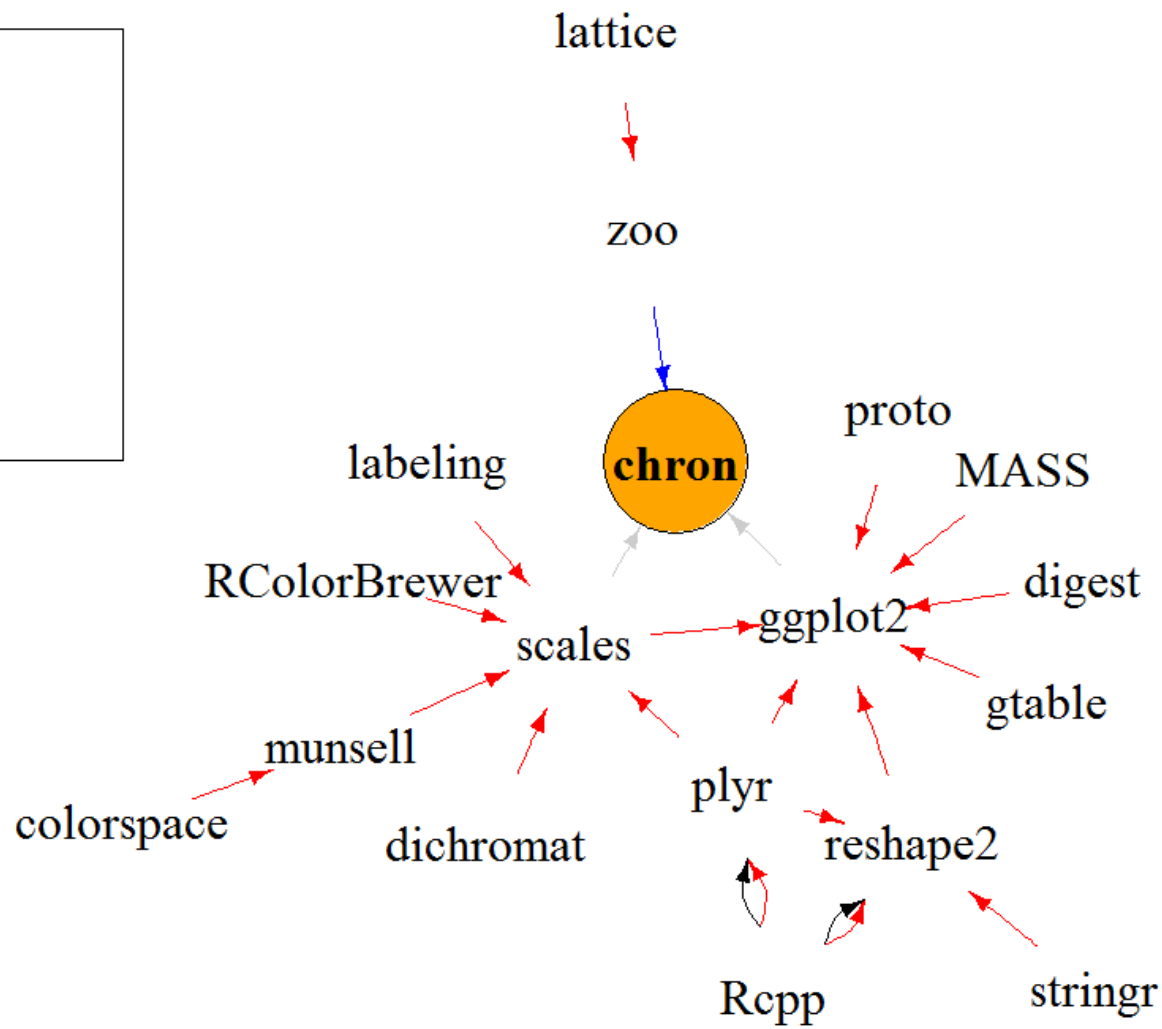# The network structure of CRAN

Andrie de Vries
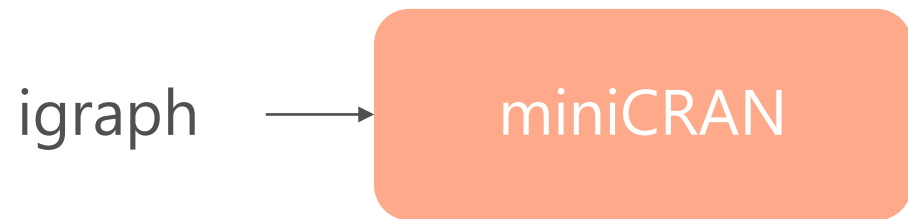
adevries@microsoft.com
@RevoAndrie

UseR!2015, Aalborg

# A network of package dependencies
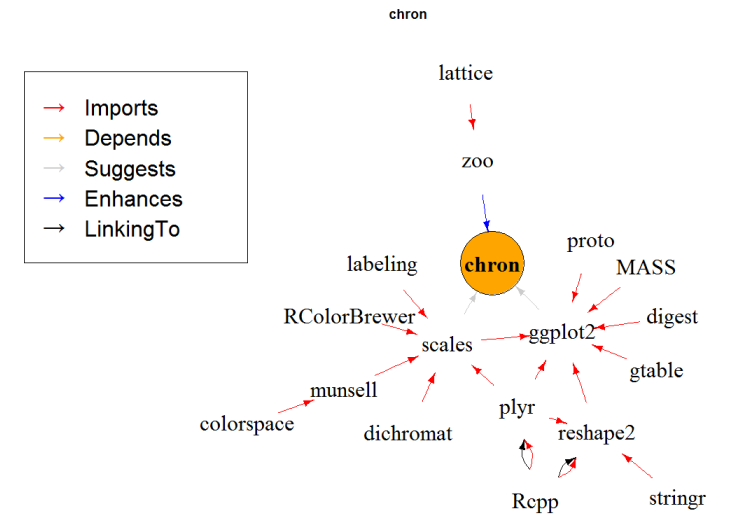
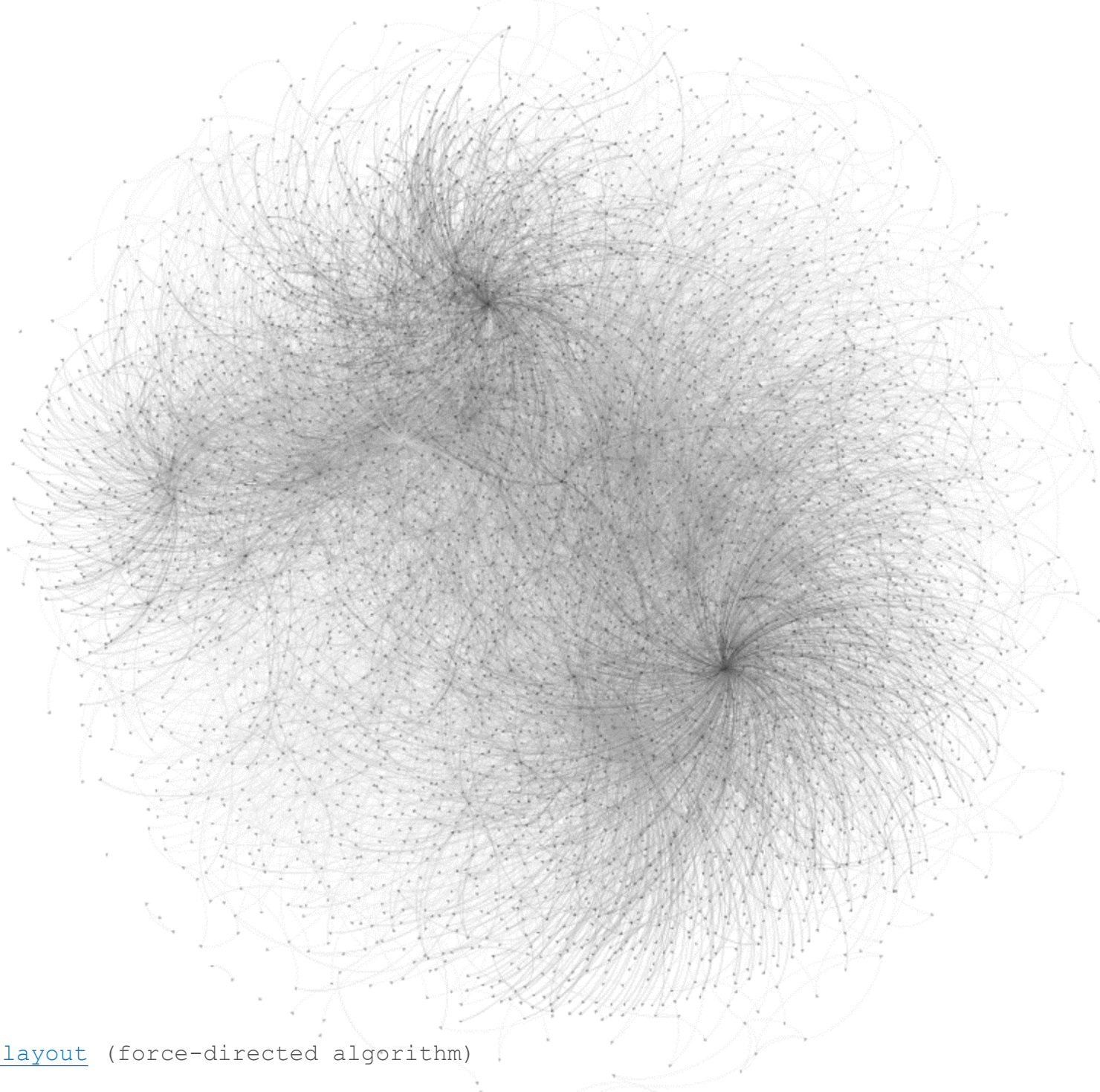# Using miniCRAN to create the network graph

igraph → **miniCRAN**



Legend:
- → Imports (red)
- → Depends (orange)
- → Suggests (grey)
- → Enhances (blue)
- → LinkingTo (black)

```
pdb <- available.packages()
tags <- "chron"
library("miniCRAN")
dg <- makeDepGraph(tags, availPkgs = pdb)
plot(dg)
```

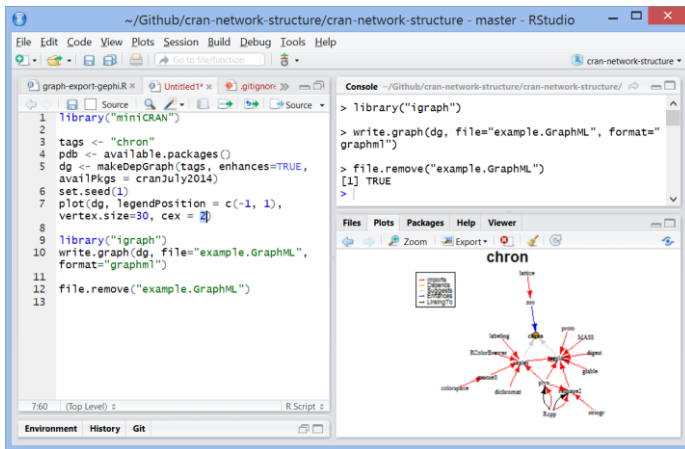# How to visualize a network of 6,791 CRAN packages?

# CRAN



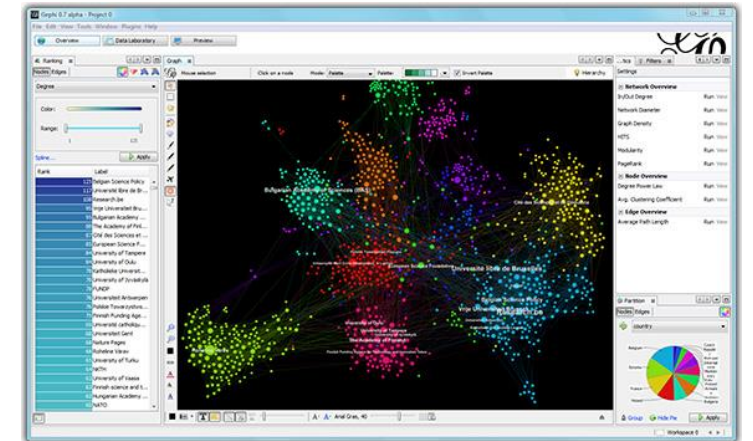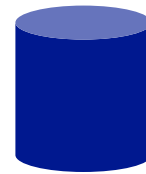Note: Fruchterman Reingold layout (force-directed algorithm)

# How to make this plot



GraphML file
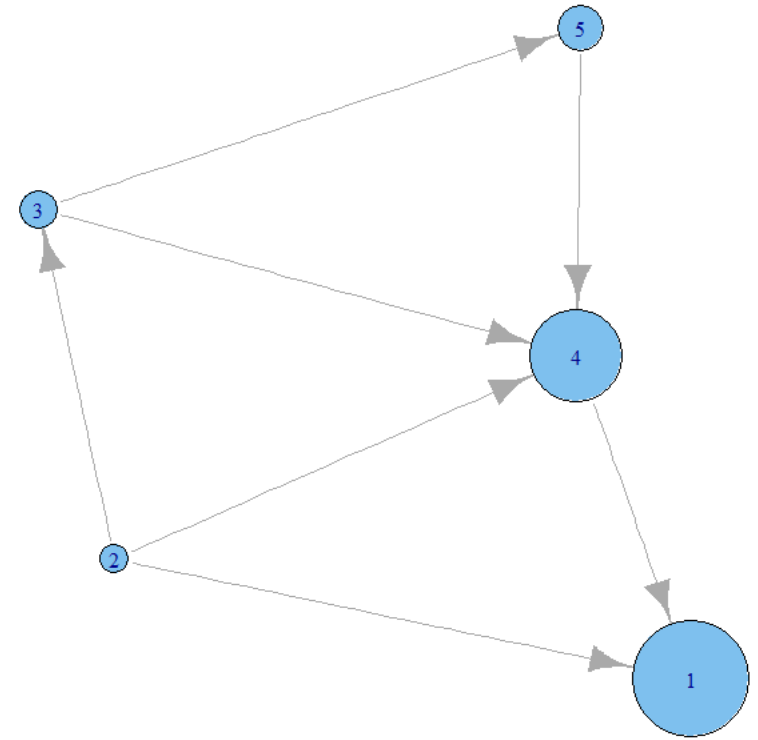
Gephi
(gephi.github.io)

```
library("igraph")
write.graph(dg,
          file = "example.GraphML",
          format = "graphml")
```

# Which nodes (packages) are more "important"?

# Reverse dependencies influence importance

- Imagine that each package dependency is a vote of confidence, cast by the community of package authors

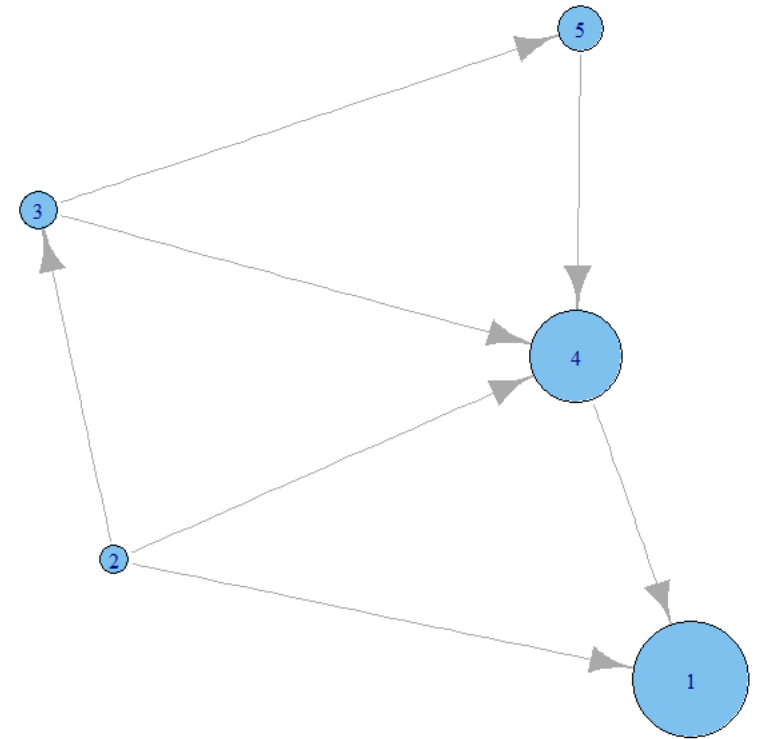- How can we compute the importance of a package to downstream packages?

# Page rank

- Page rank is famous as one of the criteria in google search results

- You can use page rank on any graph to compute "importance"
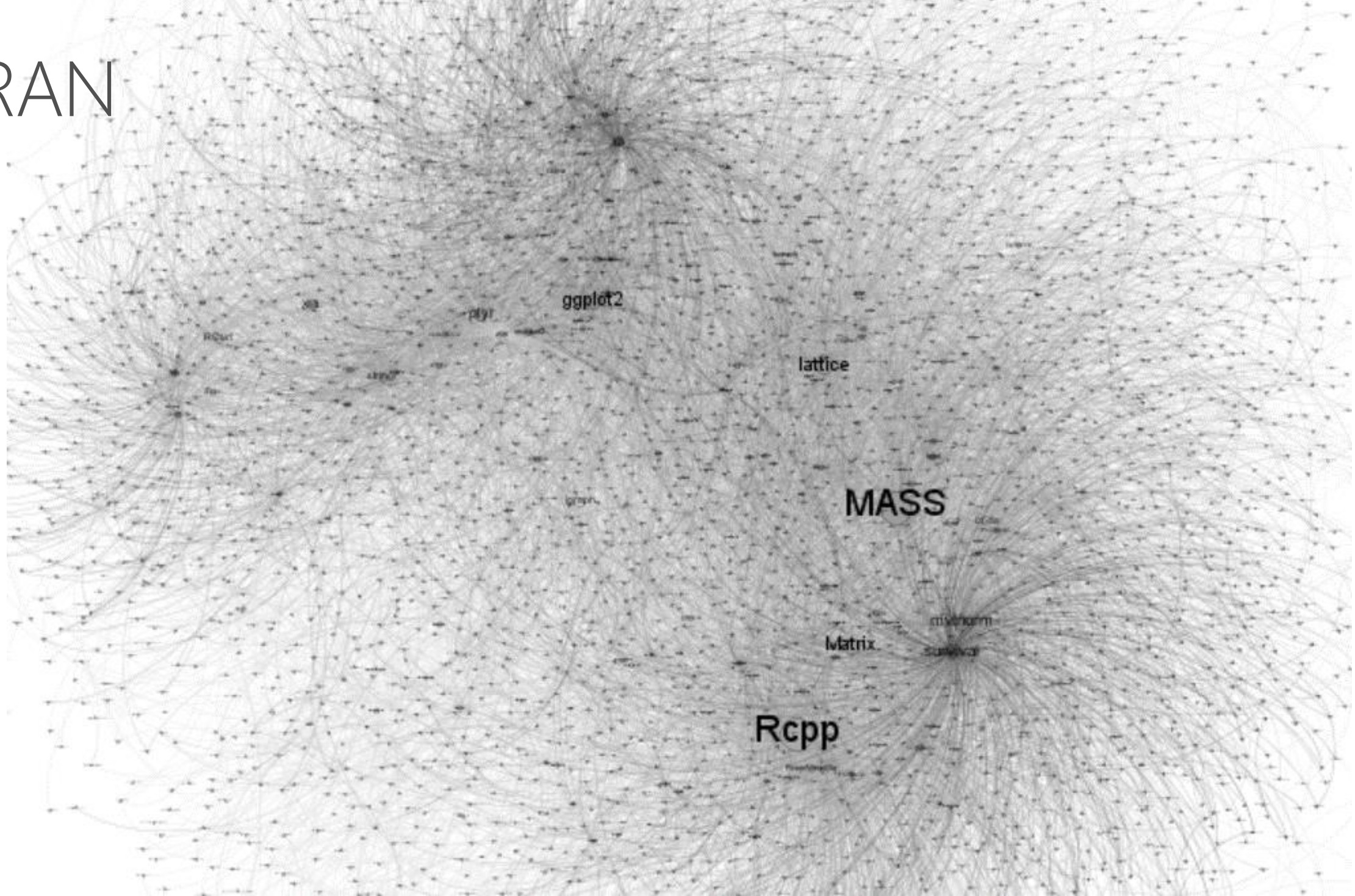
```
igraph::page.rank()
```



Simple graph with vertex size proportional to page.rank

# Pagerank results

| Package | PageRank (Dec 2014) | PageRank (Jun 2015) | Description |
|---|---|---|---|
| Rcpp | 0.0166 | 0.0197 | Interface to use C++ code in R |
| MASS | 0.021 | 0.0196 | Functions and datasets to support Venables and Ripley, 'Modern Applied Statistics with S' (4th edition, 2002). |
| Matrix | 0.01 | 0.0095 | Sparse matrix engine |
| ggplot2 | 0.0073 | 0.0086 | Graphics engine |
| lattice | 0.0096 | 0.0085 | Base R package for lattice (trellis) graphics |
| mvtnorm | 0.0088 | 0.0083 | Multivariate normal distributions |
| survival | 0.0083 | 0.0079 | Time-to-event analysis |
| plyr | 0.0067 | 0.0072 | Group-by operations |
| igraph | 0.0047 | 0.0049 | Analyse graph structures |
| XML | 0.0047 | 0.0047 | Parse and manipulate documents in XML format |

http://blog.revolutionanalytics.com/2014/12/a-reproducible-r-example-finding-the-most-popular-packages-using-the-pagerank-algorithm.html
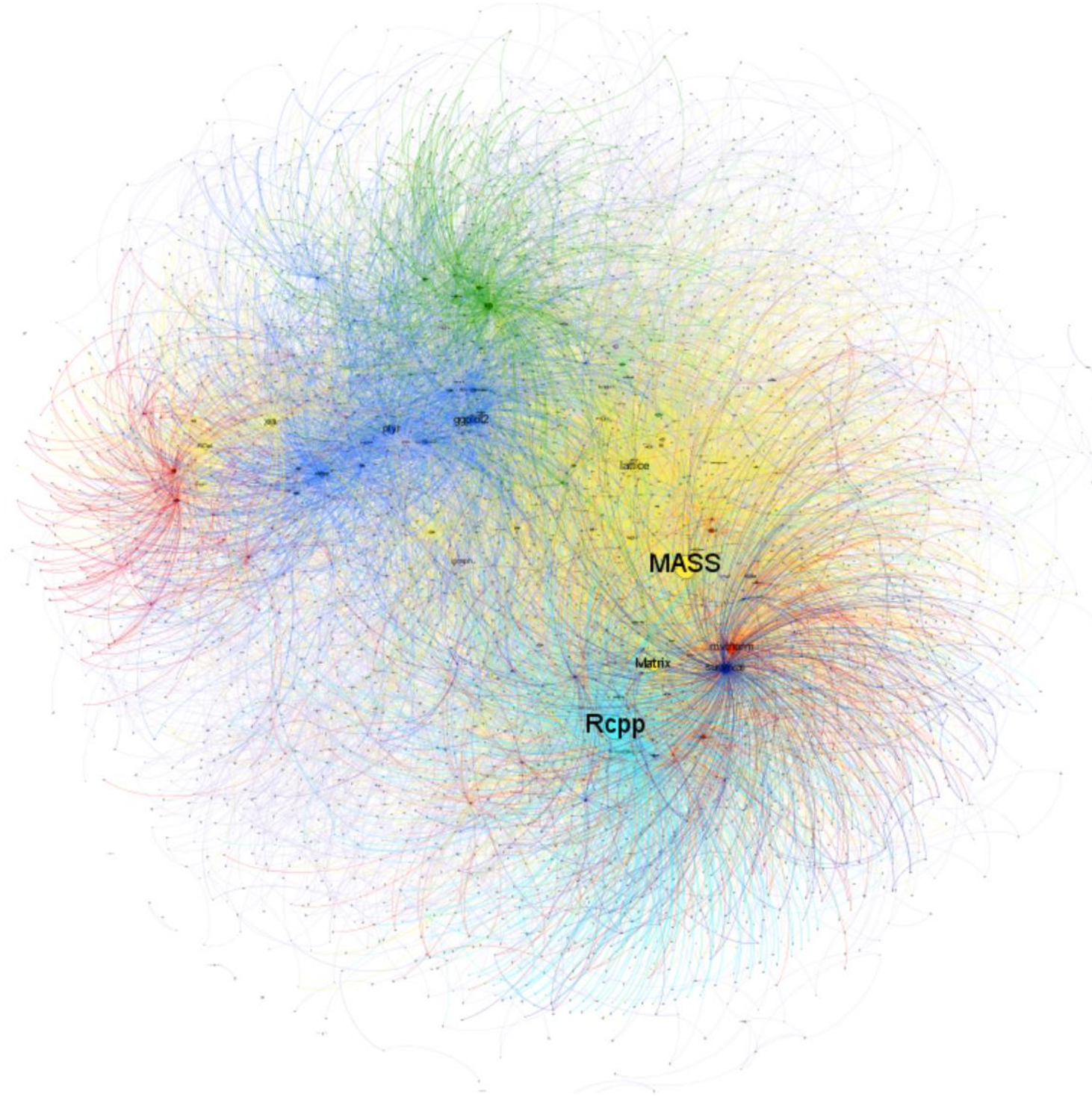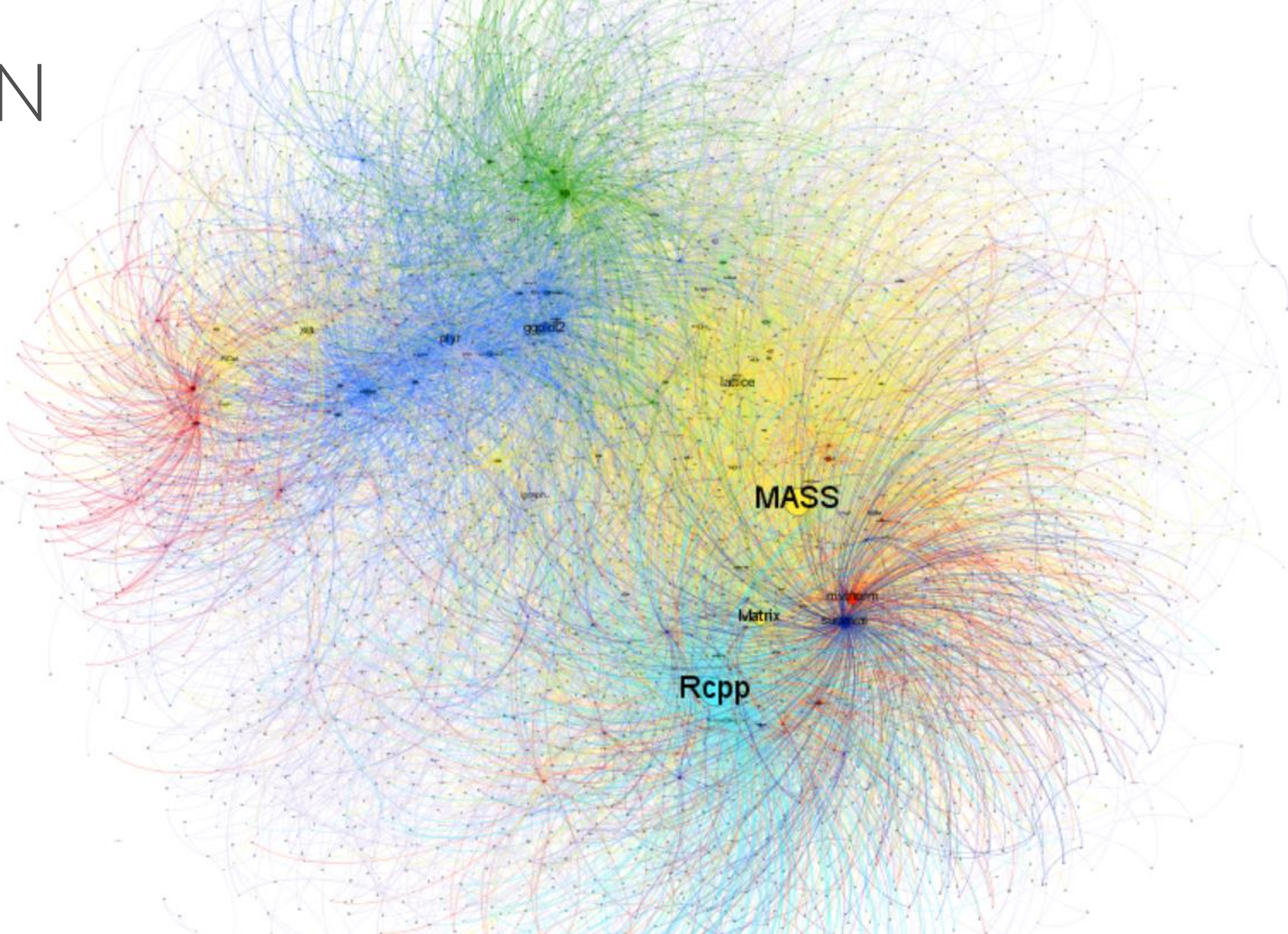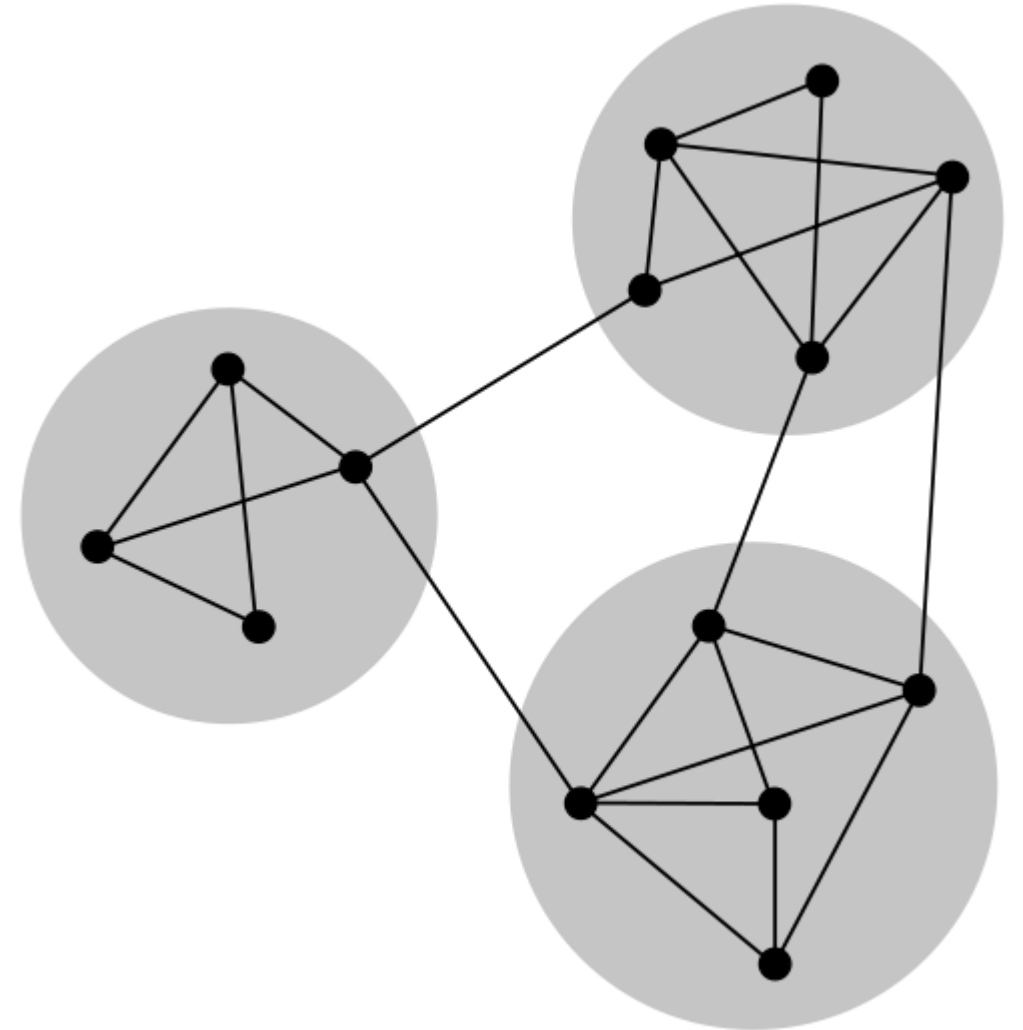
CRAN

# Can we detect clusters of packages?
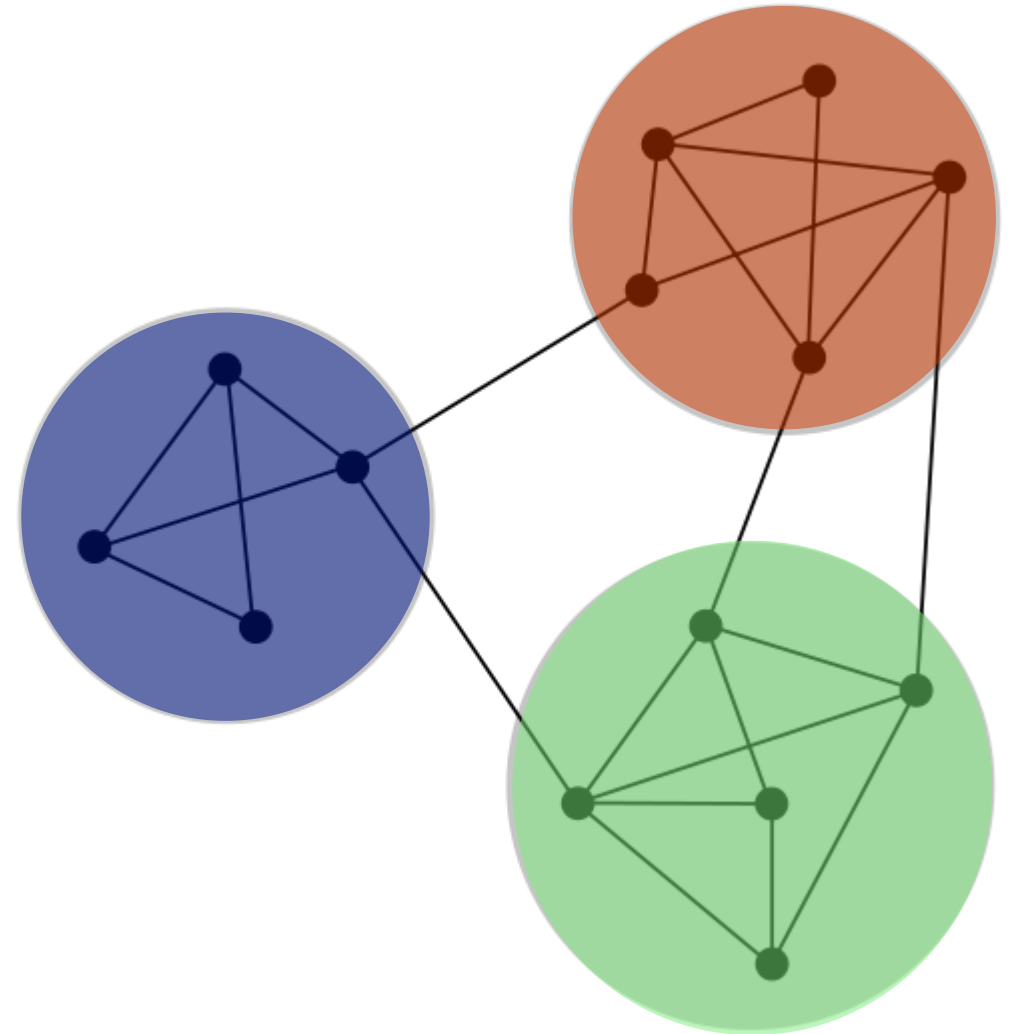
# CRAN

# CRAN

# Community detection

- Communities are densely connected groups of nodes in a graph

- Several implemented algorithms in igraph:
  - Fast greedy
  - Walktrap
  - Spinglass
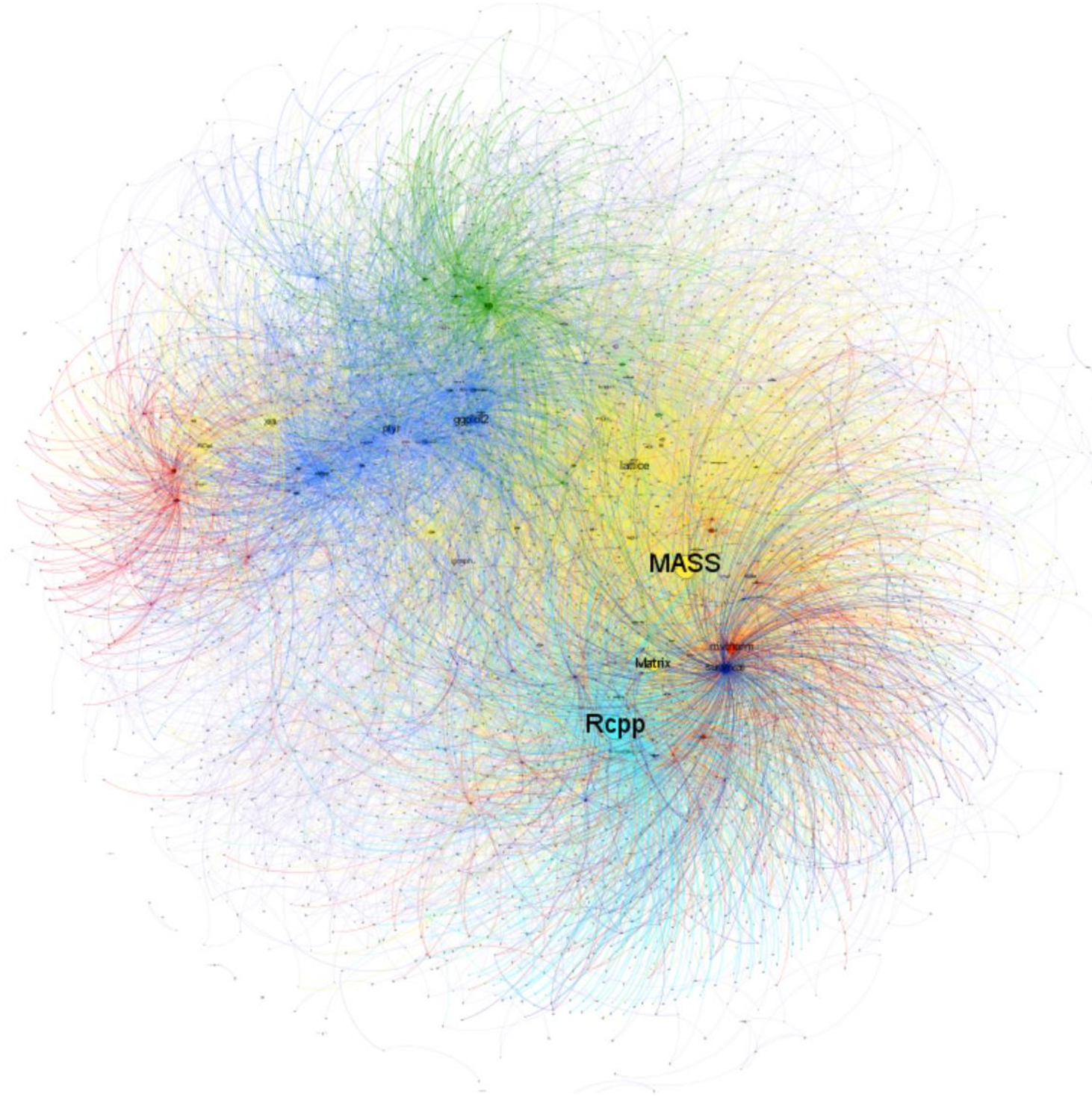  - Leading eigenvector
  - Edge betweenness

# The walktrap algorithm

- `?walktrap.community`

  - This function tries to find densely connected subgraphs, also called communities in a graph via random walks. The idea is that short random walks tend to stay in the same community.
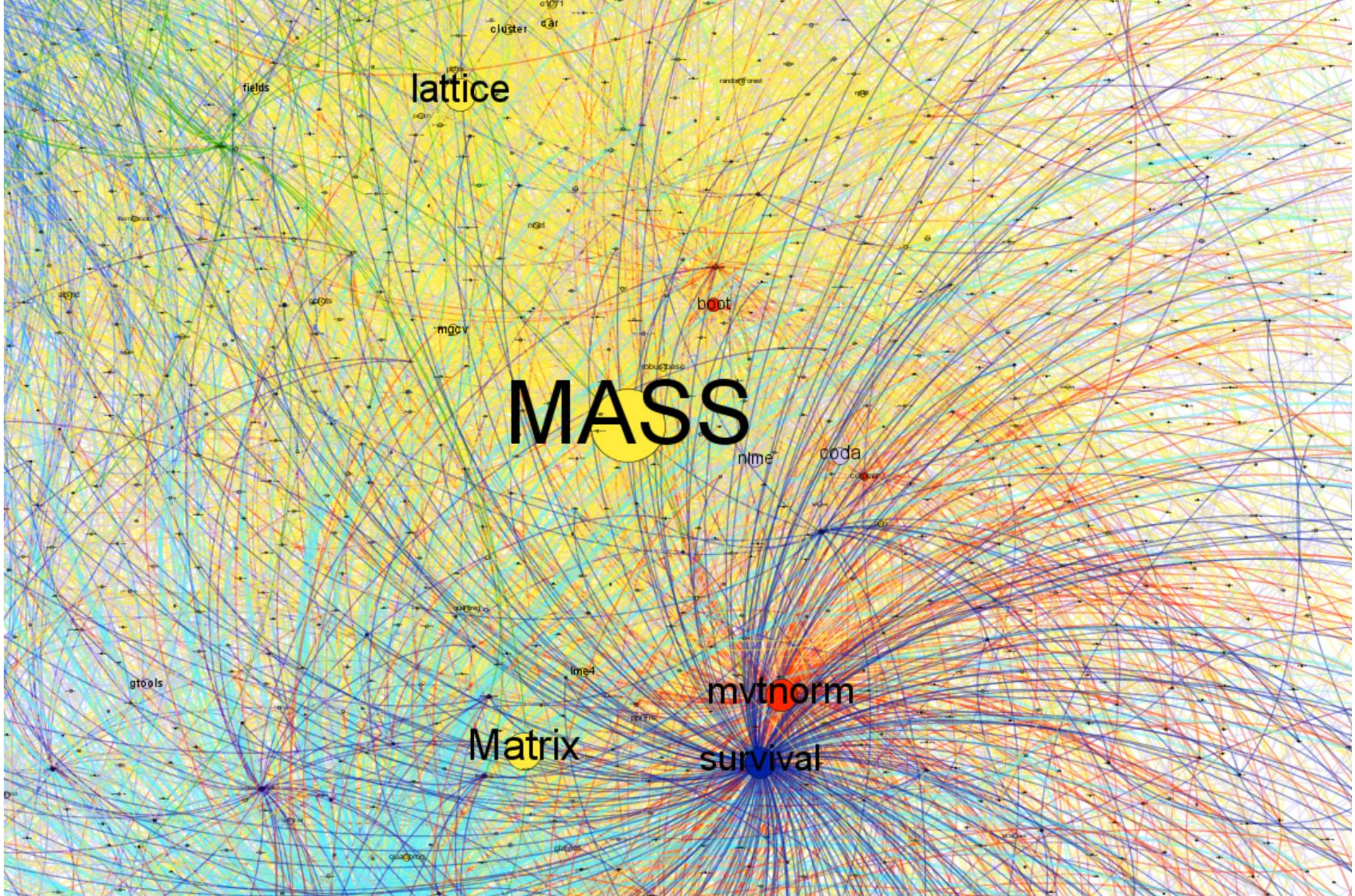
- Efficient algorithm for large, densely connected graphs*

# CRAN

# Zooming in

# Is BioConductor different?

# BioConductor



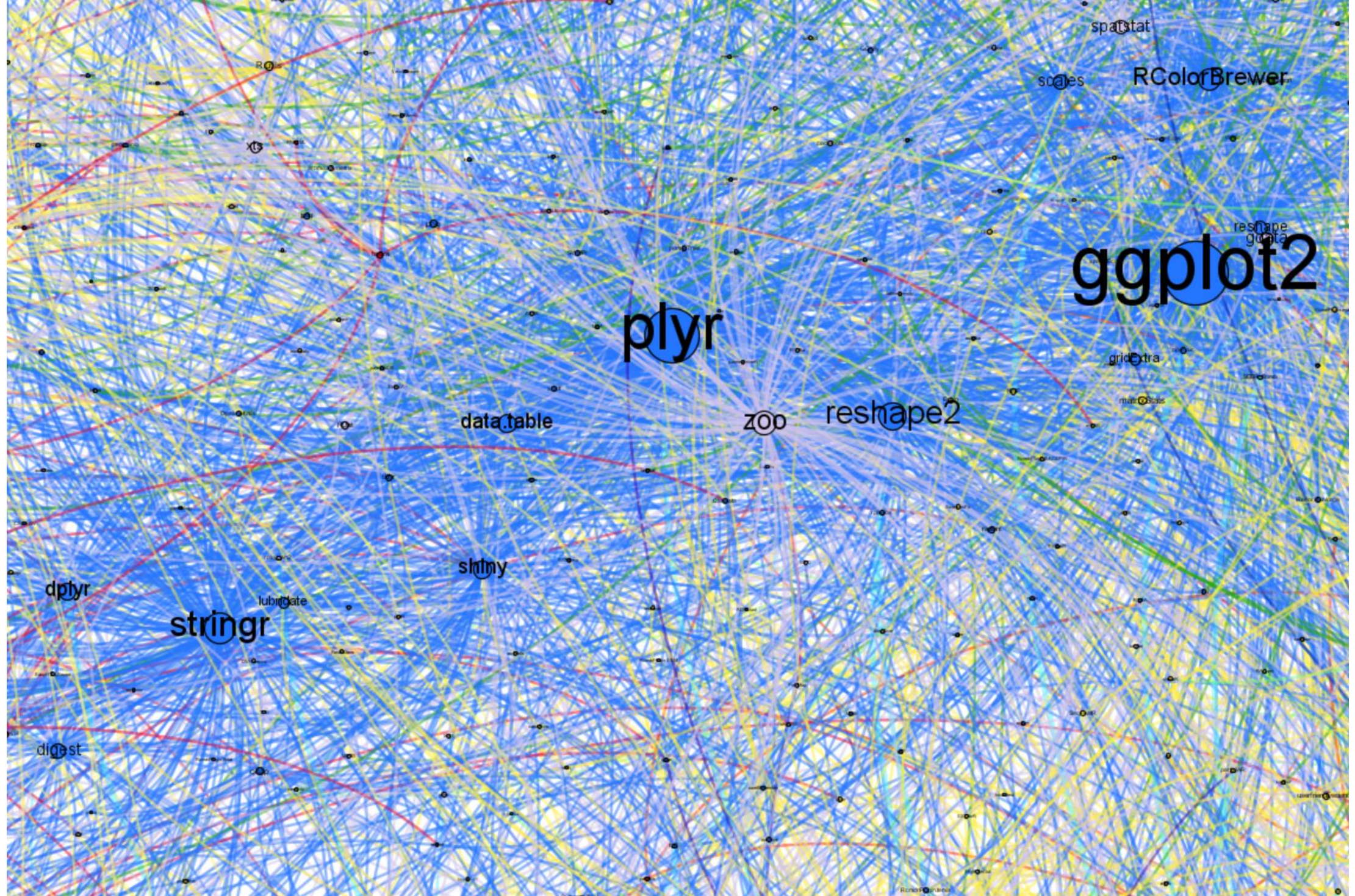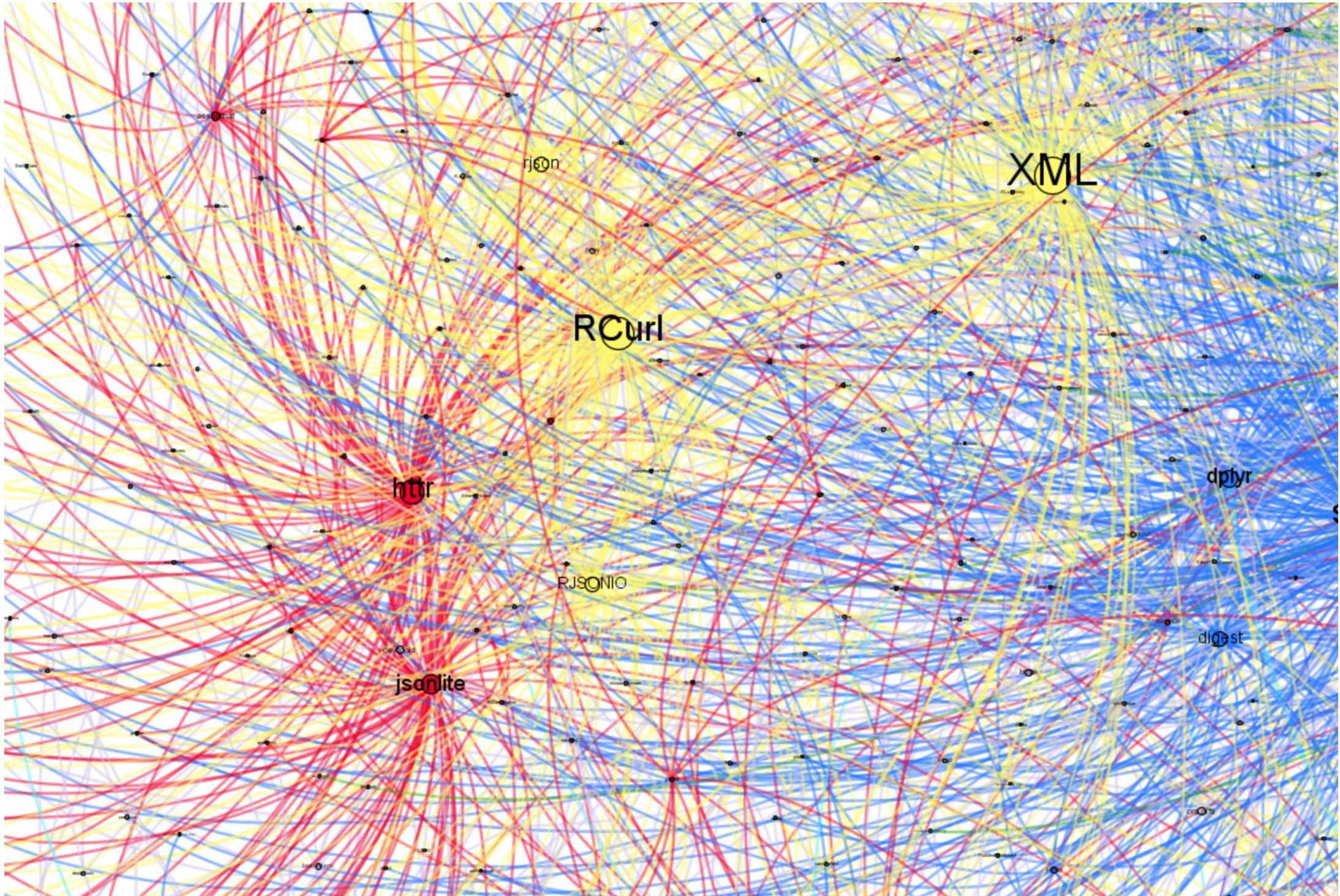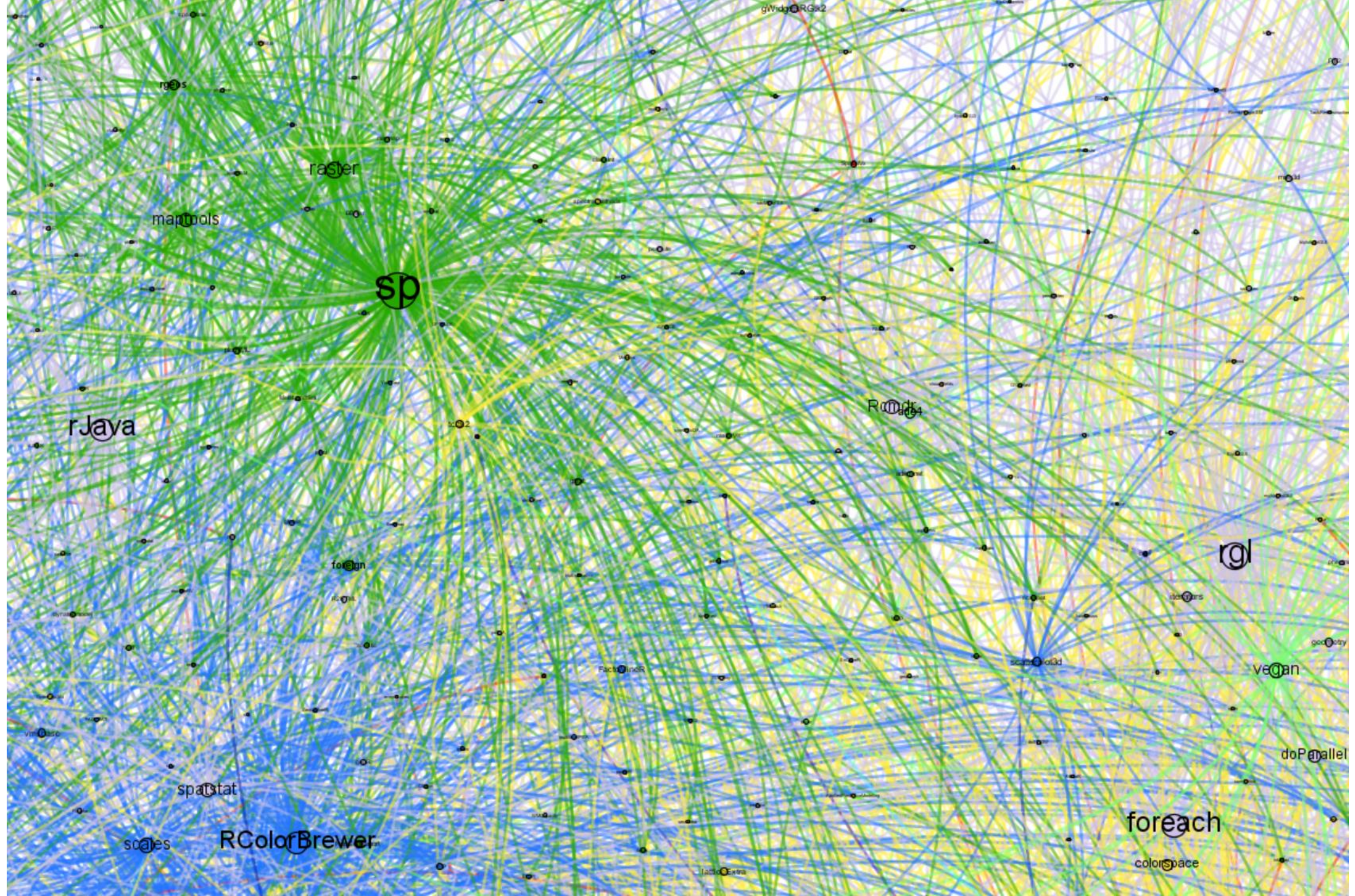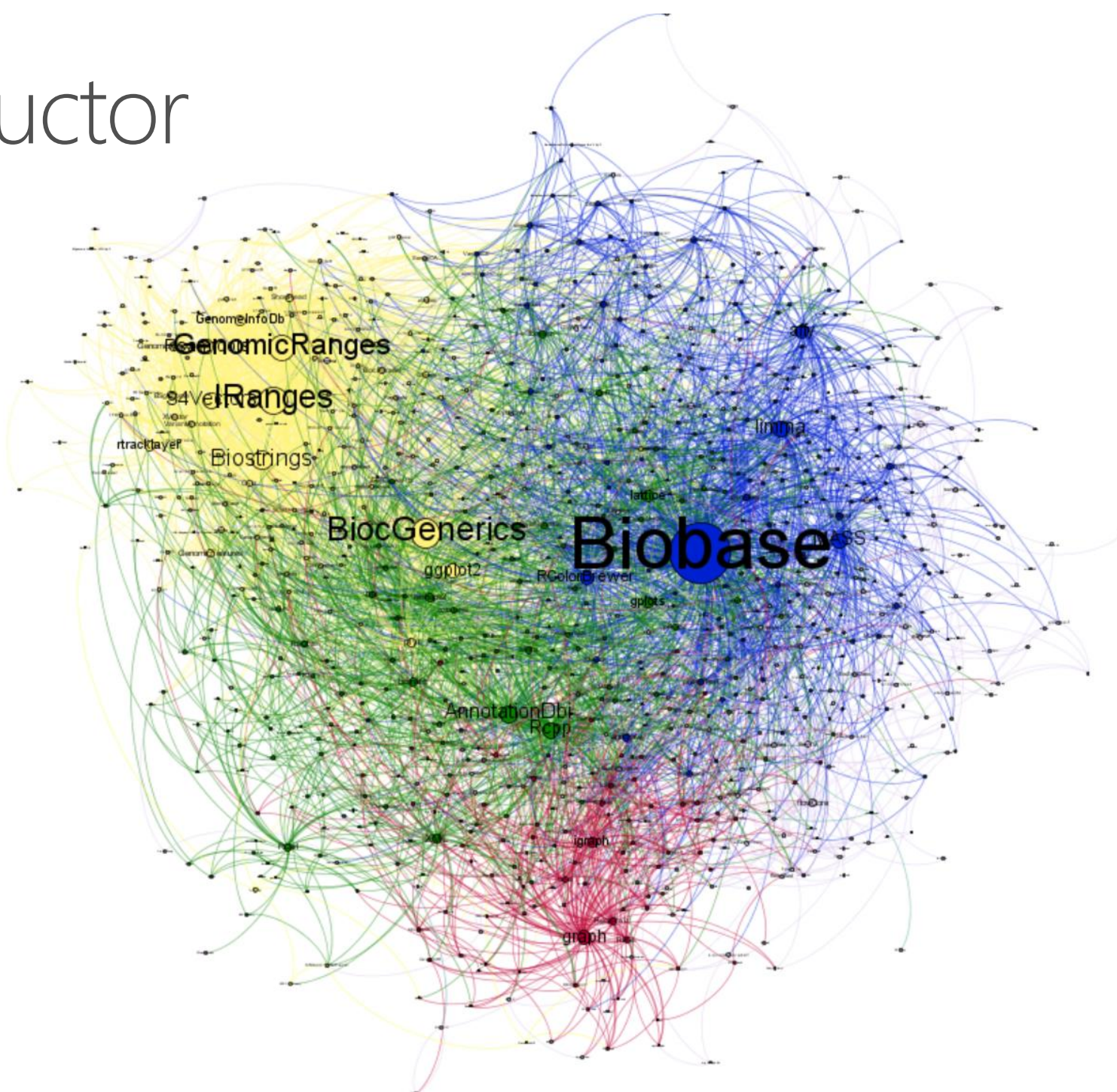Further analysis required

# Summary

# Summary

```
library(miniCRAN)
makeDepGraph()

library(igraph)
page.rank()
walktrap.community()
write.graph()
```

Scripts available at:

https://github.com/andrie/cran-network-structure

adevries@microsoft.com

@RevoAndrie